

# Image Processing and Computer Vision

Prof. Giuseppe Lisanti  
[giuseppe.lisanti@unibo.it](mailto:giuseppe.lisanti@unibo.it)

# Images...

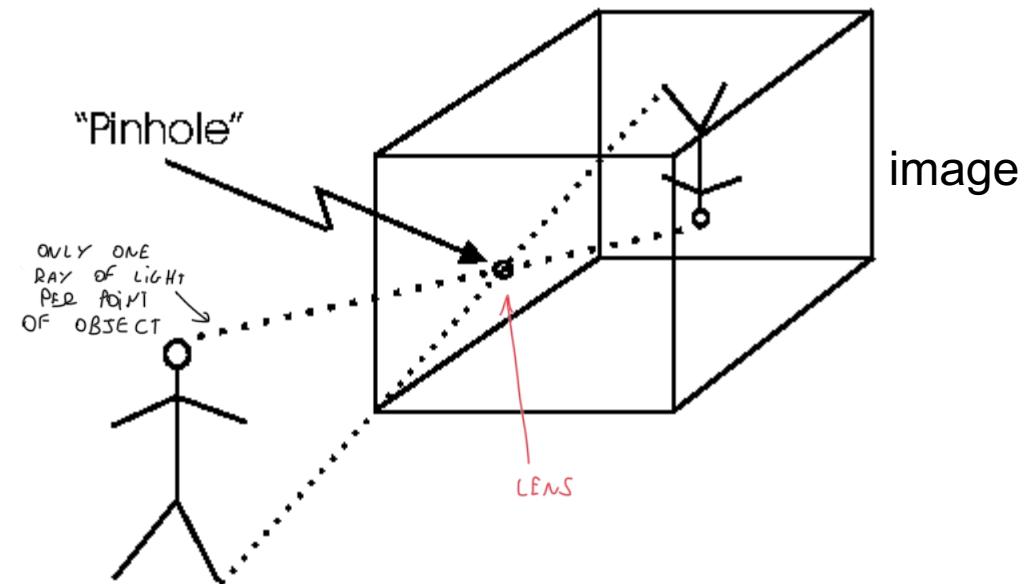
---

- An imaging device gathers the light reflected by 3D objects to create a 2D representation of the scene (i.e. the image)
- In **computer vision** we basically try to invert such a process, so as to infer knowledge on the objects from one or more digital images
- Image formation and acquisition process:
  - The geometric relationship between scene points and image points
  - The radiometric relationship between the brightness of image points and the light emitted by scene points
  - The image digitization process

# Pinhole camera model

---

- The “pinhole camera” is the simplest imaging device: light goes through the very small pinhole and hits the image plane
- Geometrically, the image is achieved by drawing straight rays from scene points through the hole up to the image plane
- Its remarkably simple geometrical model turns out to be a good approximation of the geometry of image formation in most modern imaging devices
  - however, useful images can hardly be captured by means of a pinhole camera



# Perspective Projection

- The geometric model of image formation in a pinhole camera is known as **Perspective Projection**

MAiuswLF → real world  
MINUSwLF → image

$M$  : scene point

$m$  : corresponding image point

$I$  : image plane

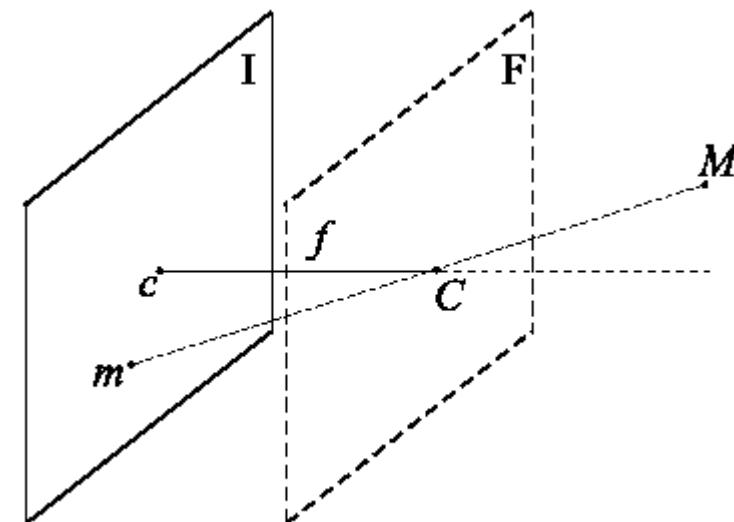
$C$  : optical centre (pin hole)

Optical axis: line through  $C$  and orthogonal to  $I$

$c$  : intersection between optical axis and image plane (image centre or piercing point)

$f$  : focal length

$F$  : focal plane



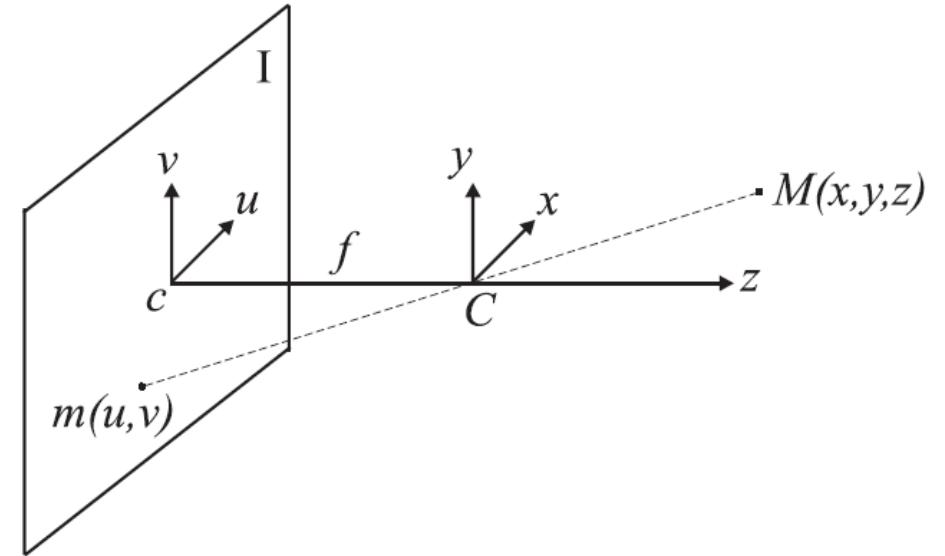
$M$  and  $m$

- We want to find a relationship between 3D and 2D points

# Perspective Projection

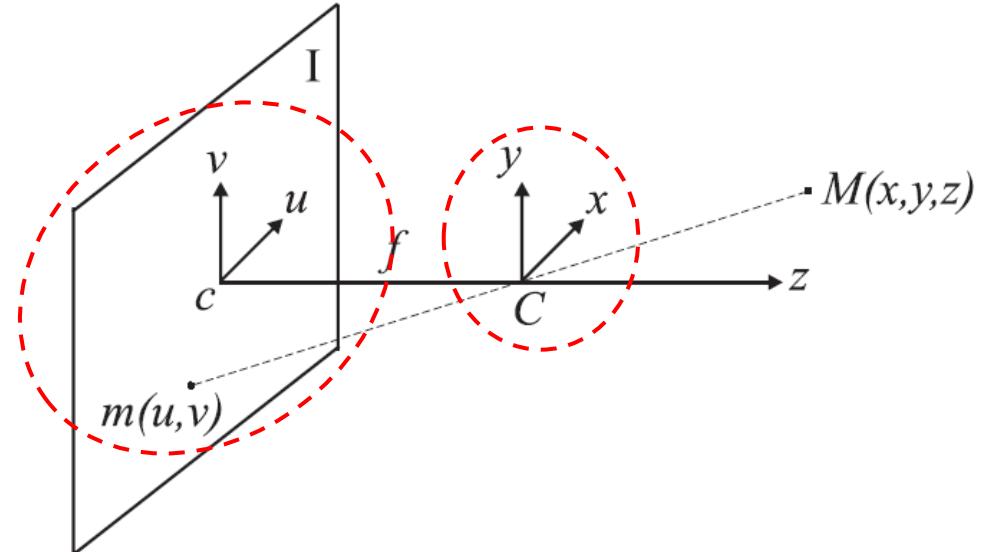
---

- Given the reference frame in figure:
  - " $u$ " is the horizontal axis in the image plane
  - " $v$ " is the vertical axis in the image plane
  - " $X$ " and " $Y$ " are the respective axis in the 3D reference system



# Perspective Projection

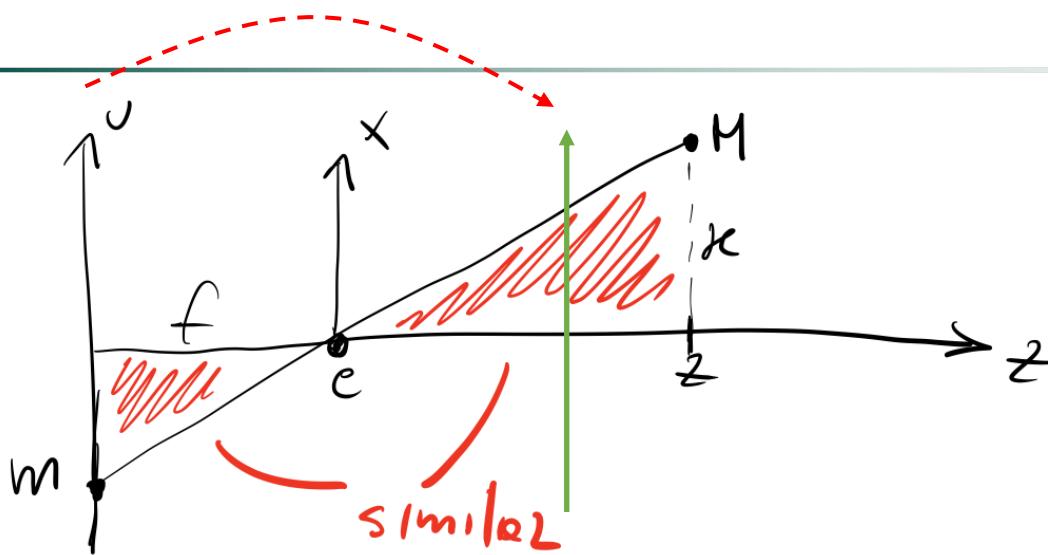
- Given the reference frame in figure:
  - " $u$ " is the horizontal axis in the image plane
  - " $v$ " is the vertical axis in the image plane
  - " $X$ " and " $Y$ " are the respective axis in the 3D reference system
    - called **camera reference system**, because it is "attached" to the camera
- For the perspective model these axis must be parallel
- The equations to map scene points into their corresponding image points are defined as follows:



$$\frac{u}{x} = -\frac{f}{z} \Rightarrow u = -x \frac{f}{z}$$

$$\frac{v}{y} = -\frac{f}{z} \Rightarrow v = -y \frac{f}{z}$$

# Perspective Projection



TRIANGLE SIMILARITY

$$\frac{u}{x} = -\frac{f}{z} \Rightarrow u = -x \frac{f}{z}$$

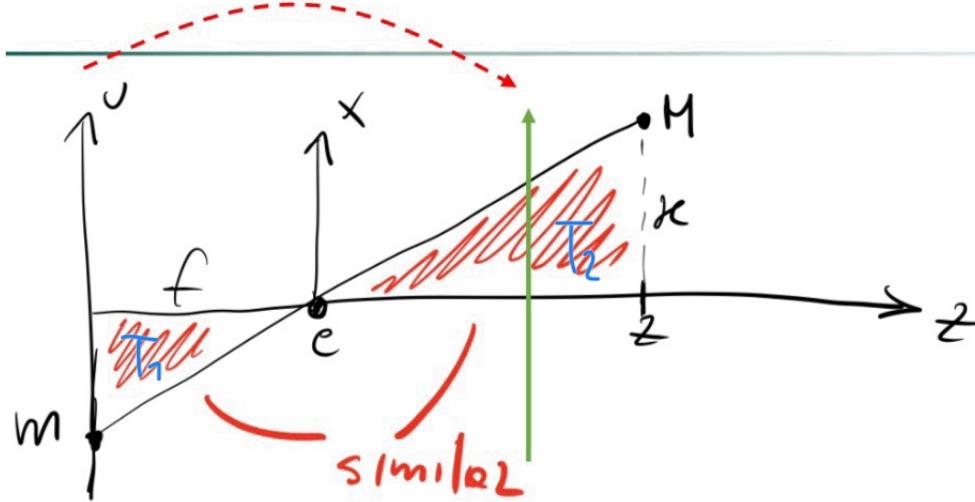
$$\frac{v}{y} = -\frac{f}{z} \Rightarrow v = -y \frac{f}{z}$$

$$\Rightarrow \frac{u}{x} = \frac{v}{y} = -\frac{f}{z}$$

- The minus means the axis get inverted
- How do we get rid of the up-down and left-right inversions?
  - Change of sign => the image plane can be thought of as lying in front rather than behind the optical centre (in real world we do not get flipped images)

$$u = x \frac{f}{z}; \quad v = y \frac{f}{z}$$

Demonstration



$T_1, T_2$  similar  $\rightarrow$  intercepting lines

$$\overline{mc} : \overline{M_c} = f : \overline{cz}$$

$c = \emptyset$  (center)

$$\frac{\overline{mc}}{\overline{M_c}} = \frac{f}{z}$$

$$\frac{\sqrt{f^2 + m_z^2}}{\sqrt{x^2 + z^2}} = \frac{f}{z}$$

... .

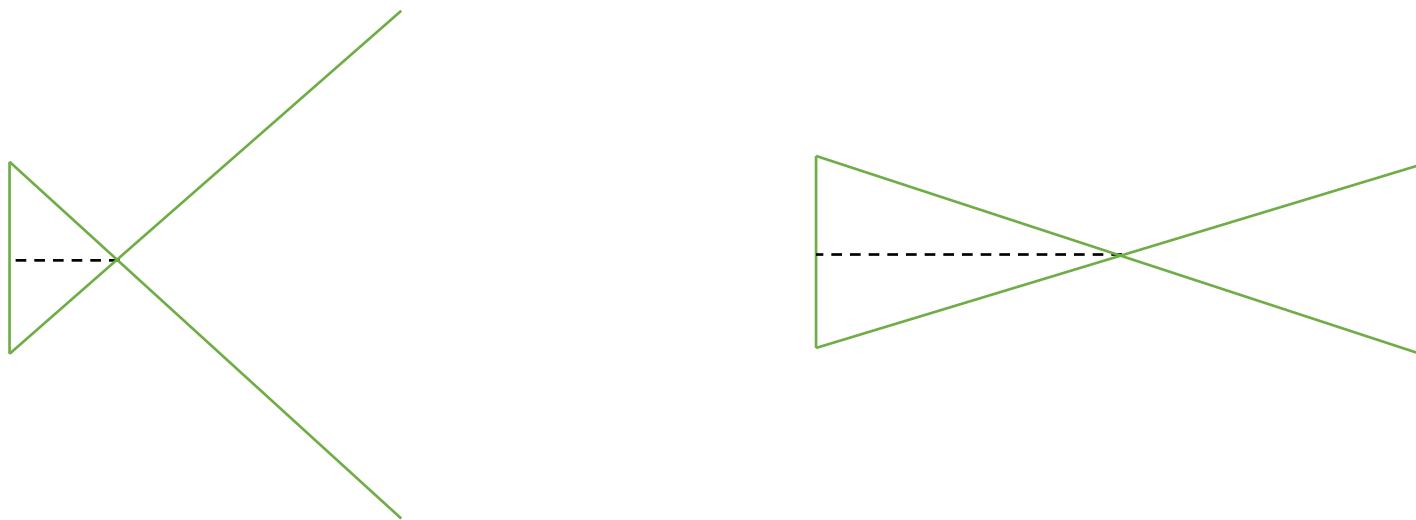
# Perspective Projection

---

- Image coordinates are a scaled version of scene coordinates (function of depth)

$$u = x \frac{f}{z}; \quad v = y \frac{f}{z}$$

- How do they scale?
  - The farther the point the smaller the coordinates (object distant from the camera)
  - The larger the focal length the bigger the object in the image (and viceversa)



# Perspective Projection

- Image coordinates are a scaled version of scene coordinates (function of depth)

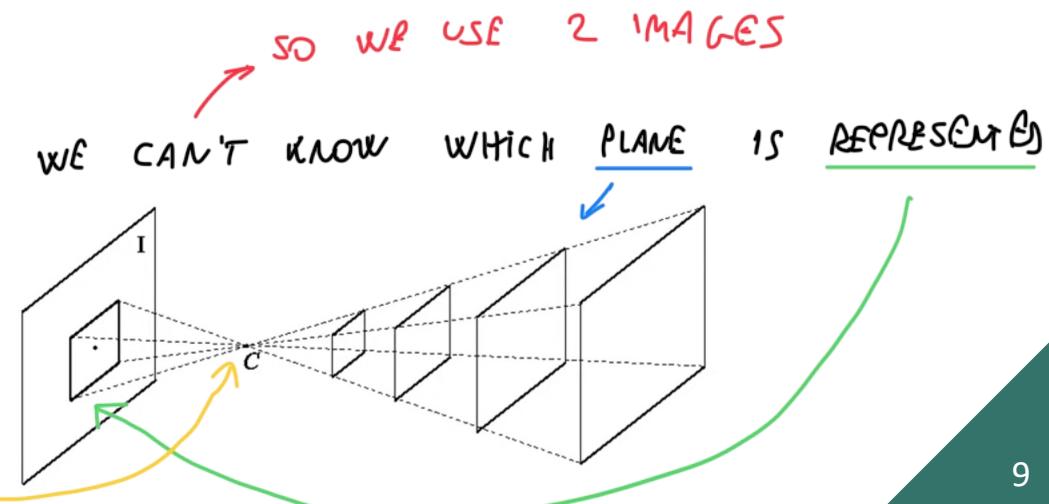
$$u = x \frac{f}{z}; \quad v = y \frac{f}{z}$$

- How do they scale?

- The farther the point the smaller the coordinates (object distant from the camera)
- The larger the focal length the bigger the object in the image (and viceversa)

- We scale the world inversely with respect to the depth

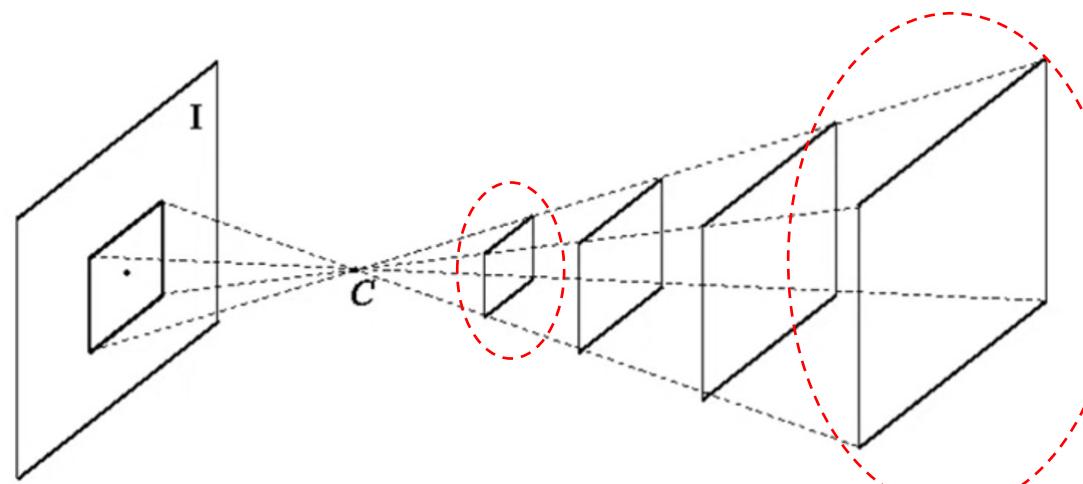
RELATION IS NOT BIJECTIVE



- The image formation process deals with mapping a 3D space onto a 2D space => loss of information

# Perspective Projection

- The mapping is not a bijection: a given **scene point** is mapped into a **unique image point**
  - a given **image point** is mapped onto a **3D line** (i.e. the line through the point,  $m$ , and the optical centre,  $C$ ).

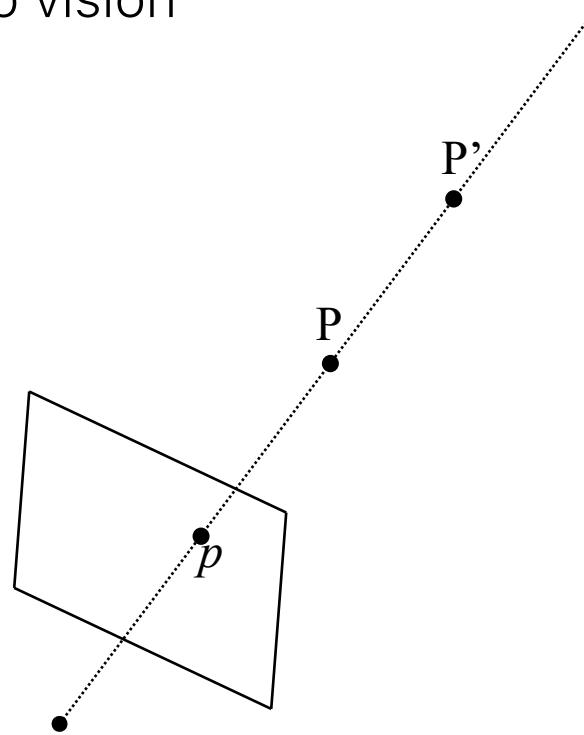


- Recovering the 3D structure of a scene from a single image is an **ill-posed problem** (the solution is not unique)
  - For an image point we can only state that its corresponding scene point lays on a line but cannot disambiguate a specific 3D point along such a line (i.e. we know nothing about the distance to the camera).
- How can we solve this?

# Stereo images allow to infer 3D

---

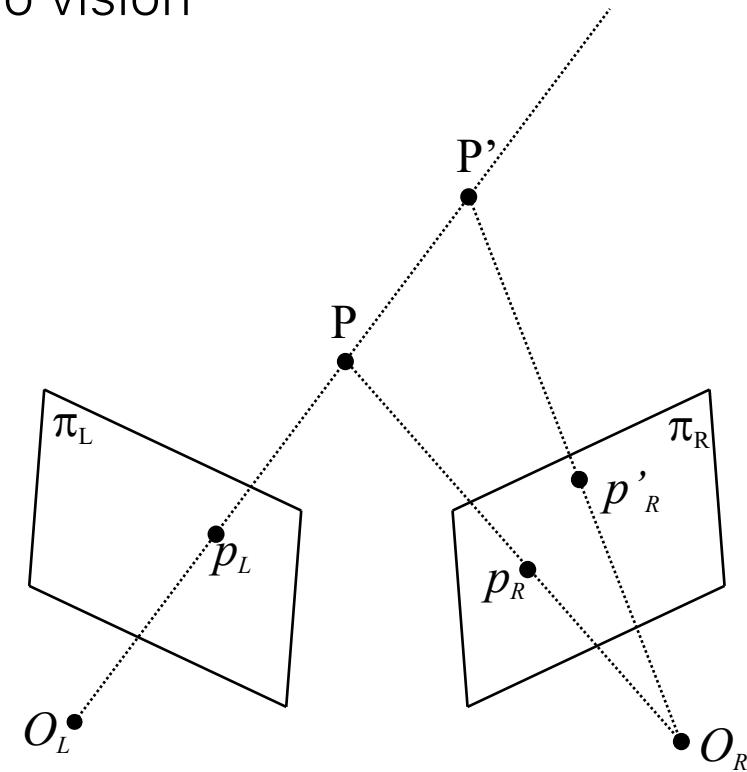
- Solution: use multiple images (at least two) => stereo vision



# Stereo images allow to infer 3D

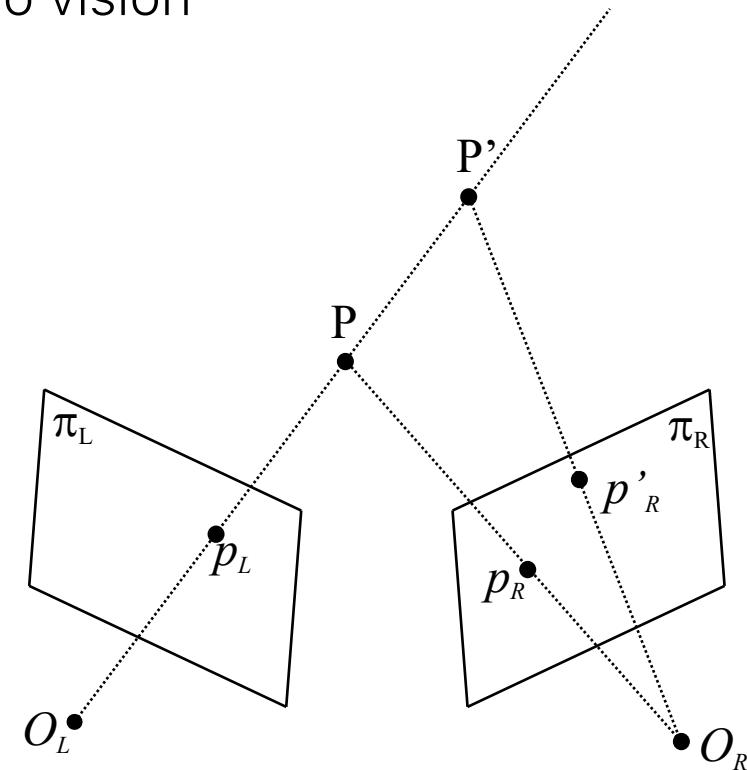
---

- Solution: use multiple images (at least two) => stereo vision
- The human visual system is a stereo vision system
- Stereo images allow to infer 3D
- Given **correspondences**, 3D information can be recovered easily by triangulation



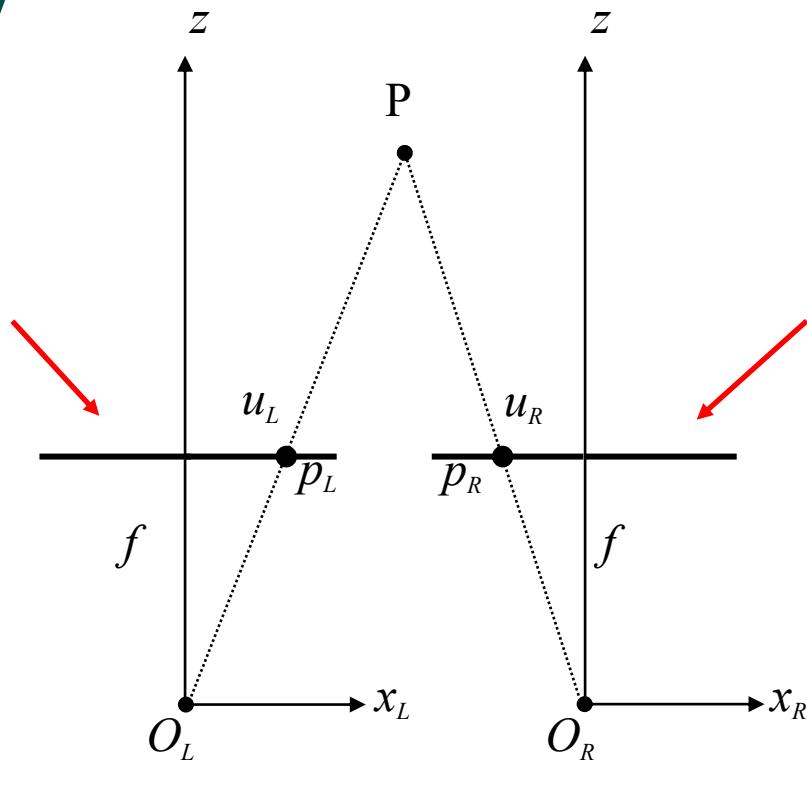
# Stereo images allow to infer 3D

- Solution: use multiple images (at least two) => stereo vision
- The human visual system is a stereo vision system
- Stereo images allow to infer 3D
- Given **correspondences**, 3D information can be recovered easily by triangulation



# Standard stereo geometry

- Assumptions:
  - Parallel  $(x, y, z)$  axes
  - Same focal length  $\Rightarrow$  coplanar image planes
- The transformation between the two reference frames is just a translation ( $b$ ), usually horizontal
- You need to **sense** two images at the very same moment
- You can put the two cameras "as you want" but they must observe the same object



$$P_L - P_R = \begin{bmatrix} b \\ 0 \\ 0 \end{bmatrix} \quad \rightarrow$$

$$\begin{aligned} x_L - x_R &= b \\ y_L &= y_R = y \\ z_L &= z_R = z \end{aligned}$$

# Standard stereo geometry

- The two cameras are displaced at a given quantity  $b$  called **baseline**  
*GIVEN*

- Disparity:** difference between the horizontal coordinates in the left and right images (horizontal displacement)

*GUARANTEES CORRESPONDENCE*

$$\begin{cases} u_L = x_L \cdot f/z \\ u_R = x_R \cdot f/z \end{cases}$$

- Inverse relation: the larger the disparity the smaller depth, and vice versa

- If a point has a large disparity it is close to the camera...

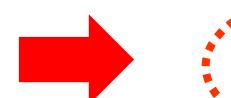
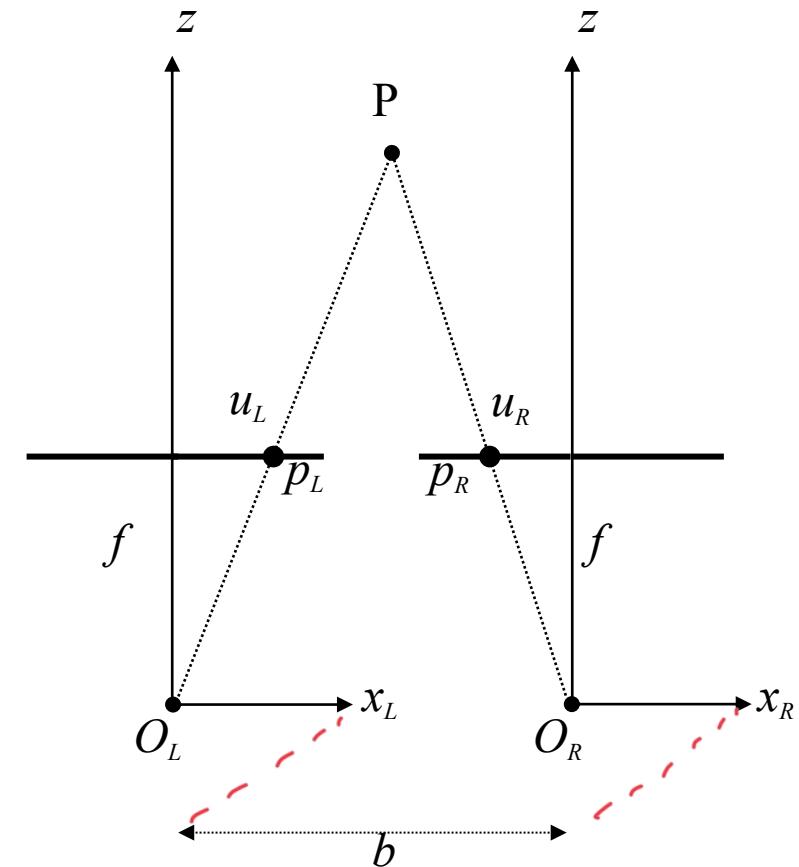
$$v_L = v_R = y \cdot f/z$$

$$u_L - u_R = b \cdot f/z$$

$$u_L - u_R = d$$

(disparity)

$$d = b \cdot f/z$$



$$z = b \cdot f/d$$

Fundamental relationship in stereo vision

# Standard stereo geometry

$$v_L = v_R = y \cdot f/z$$

- Since we are given just two 2D images there is no info about the correspondence between two points in the two images

- Given  $p_L$ , I want to find  $p_R$

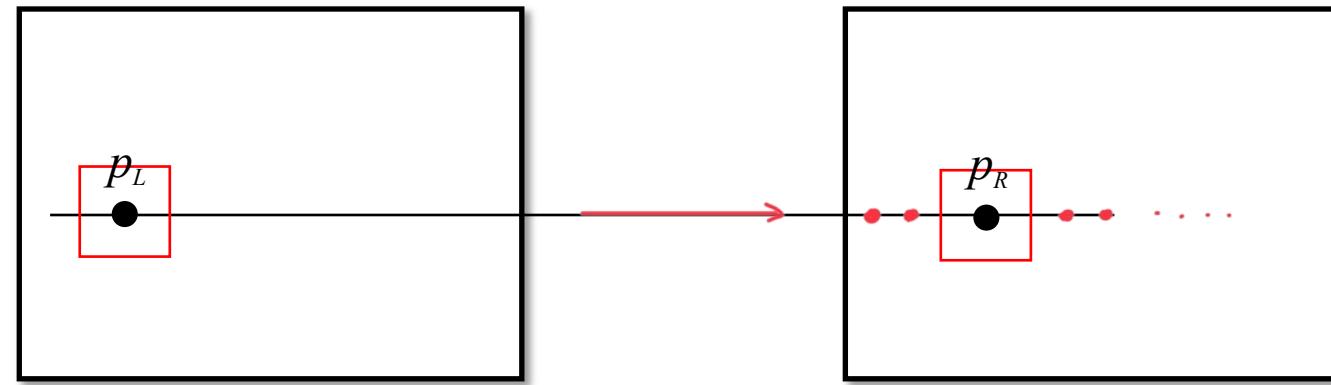
- We can search along horizontal lines

- Stereo matching

IMAGES ARE PERFECTLY ALIGNED

SO, IN ORDER TO FIND  $p_R$ , WE START FROM  $p_L$  AND  
LOOK ALL POINTS IN THE LINE ALONG THE SECOND  
IMAGE, SEEKING FOR THE MOST SIMILAR POINT

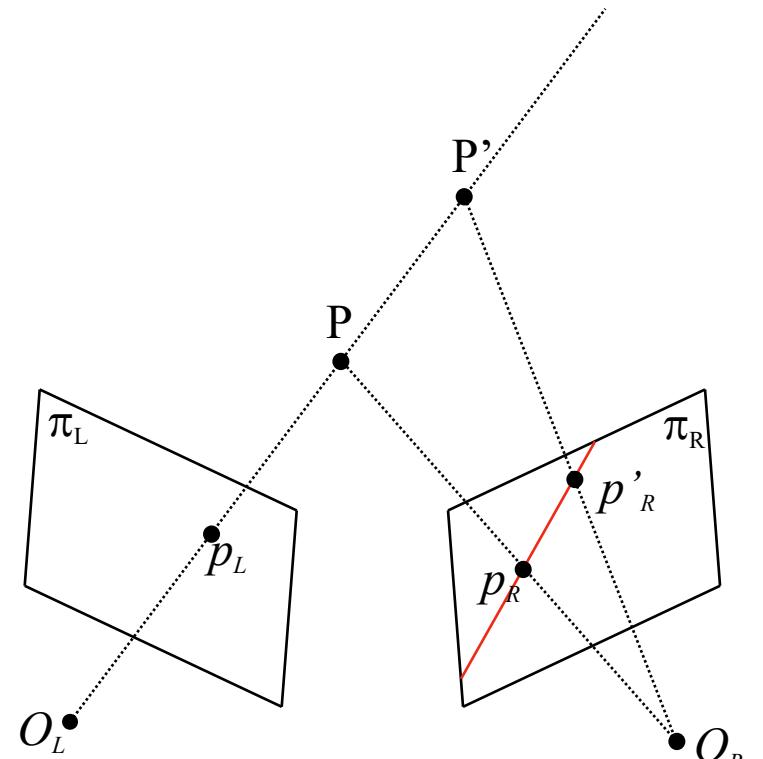
??



- For example you search for same color/pattern along the same line (better: use a window around the points across the line, block matching)

# Epipolar Geometry

- What if the two cameras are no longer aligned? Do we need to search through the whole image?
  - We can project the line related to point  $P_L$  in the right plane and search across that line
    - The search space of the stereo correspondence problem is always 1D !
- Issue: this projection can be computed only if the transformation between the two cameras is known (relative mapping between the two cameras)
  - A roto-translation and the focal lengths

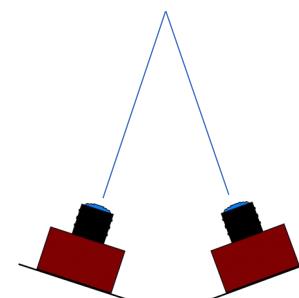
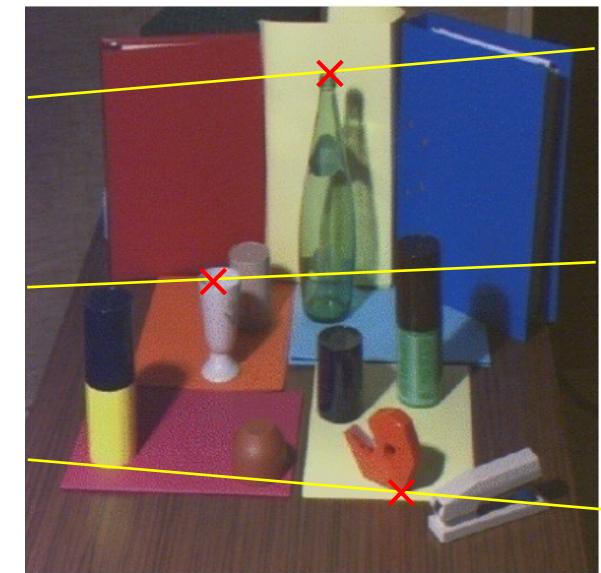
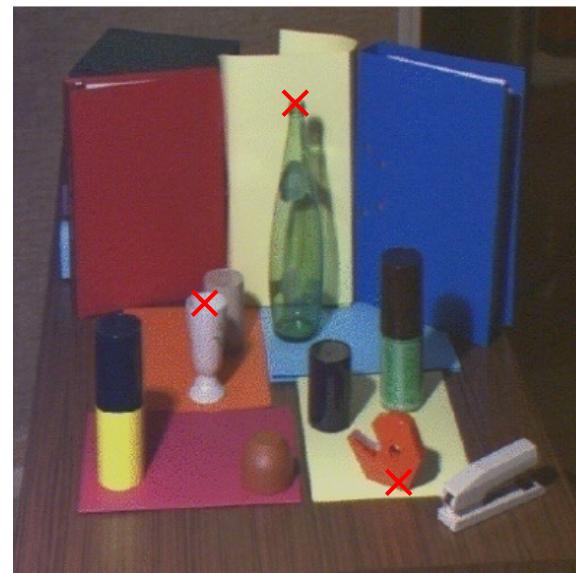


Epipolar line  
(associated with  $p_L$  in  $\pi_R$ )

# Epipolar Geometry

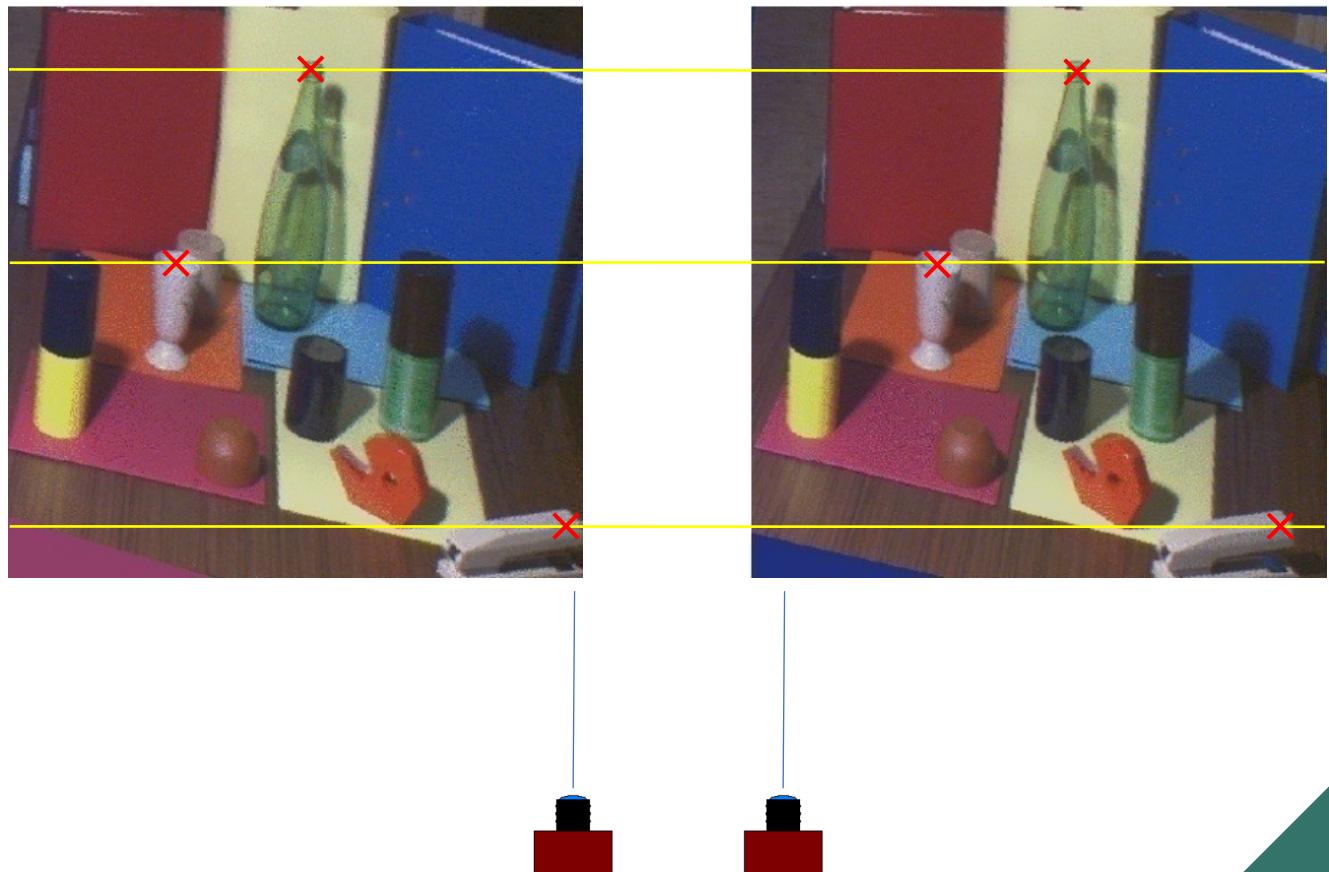
- It is almost impossible to build a stereo rig which is perfectly aligned horizontally
- Searching through oblique epipolar lines is awkward!
  - It is also computationally less efficient
- What can we do?

WE RECTIFY IT



# Rectification

- What people do in practice is to convert epipolar geometry to standard geometry (**Rectification / Warping**)
- Warp the images as if they were acquired through a standard geometry (horizontal and collinear conjugate epipolar lines)
  - Compute and apply to both images a transformation (i.e. **homography**) known as **rectification**



# Stereo Correspondence

- Given a point in one image (e.g. L) find that in the other image (R) which is the projection of the same 3D point. Such image points are called corresponding points.



$$D_1 < D_2 \quad D_1 > D_2$$

Corresponding points look similar in the two images

A window (group of points)  
could be confused with other  
windows

Points farther away have a smaller disparity, while close points have a larger disparity

# Properties of Perspective Projection

---

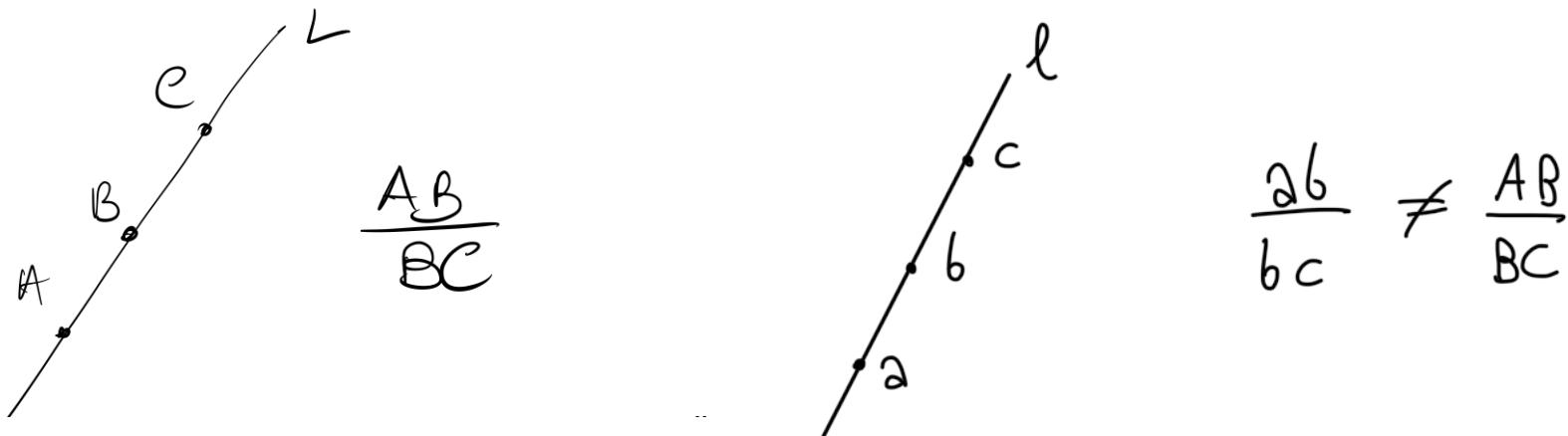
- The image of a 3D line segment of length  $L$  lying in a plane parallel to the image plane at distance  $z$  from the optical centre will exhibit a length given by:

$$l = L \frac{f}{z}$$

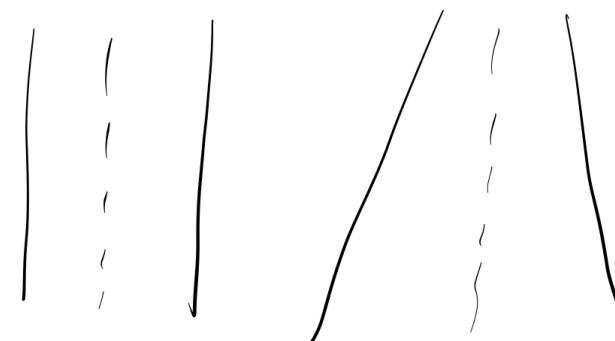
- This relationship is more complicated for an arbitrarily oriented 3D segment, as its position and orientation need to be accounted for as well
- Nonetheless, for given position and orientation, length always shrinks alongside distance

# Properties of Perspective Projection

- Perspective projection maps 3D lines into image lines
  - Ratios of lengths are not preserved (unless the scene is planar and parallel to the image plane).

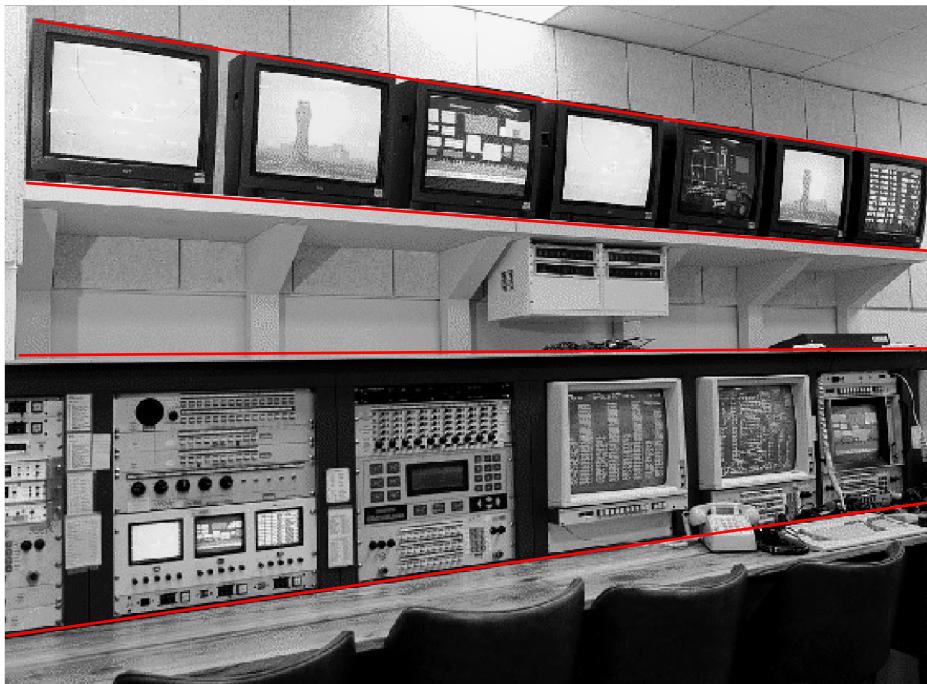


- Parallelism between 3D lines is not preserved (except for lines parallel to the image plane)

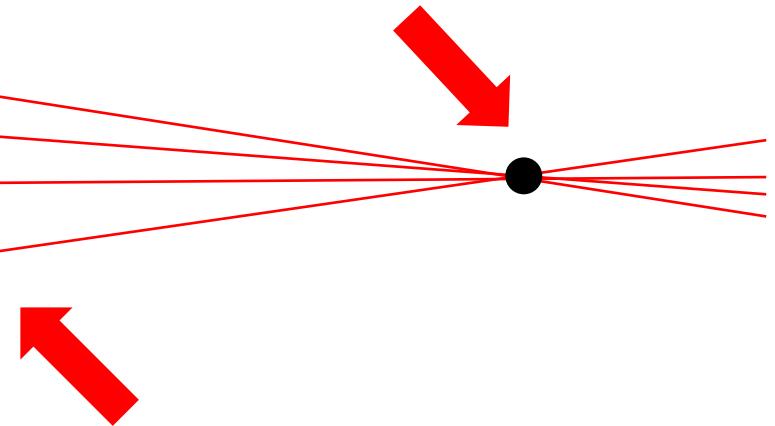


# Vanishing point

- The images of parallel 3D lines intersect at a point



Vanishing point  
(not necessarily within the image)

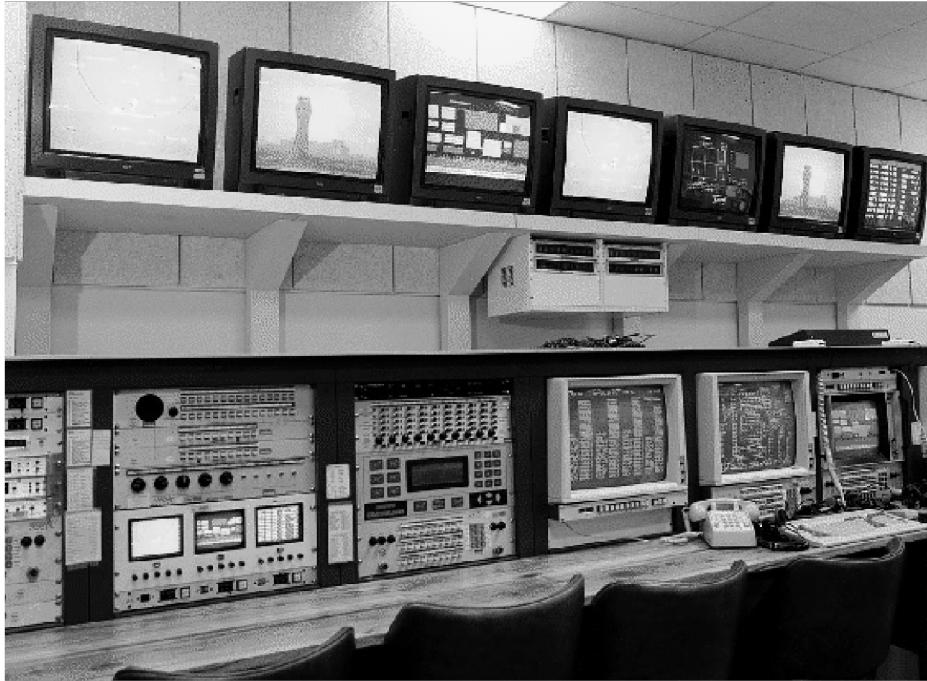


Parallel in real world!

# Vanishing point

---

- The images of parallel 3D lines intersect at a point, which is referred to as vanishing point.



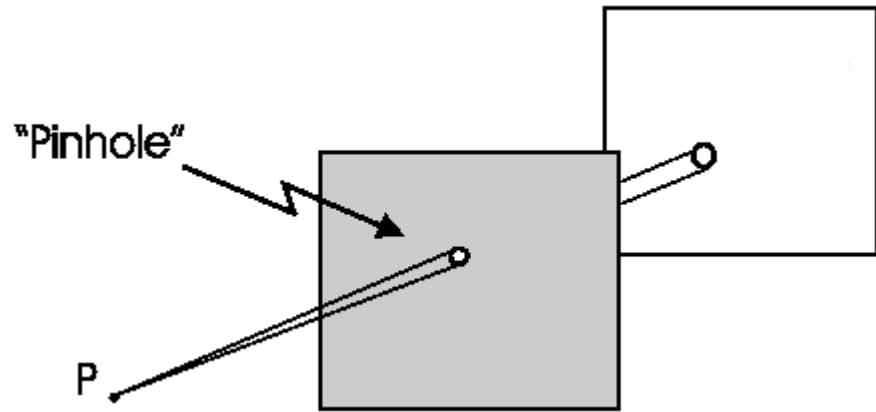
- If the lines are parallel to the image plane they meet at infinity

# Depth of Field (DOF)

WHILE TAKING A PHOTO, IF YOU STAY STILL,  
IT STARTS TO FOCUS (collecting light rays),  
BUT IF YOU MOVE YOU CHANGE THE  
GEOMETRICAL RELATIONS AND IT UNFOCUS

- A scene point is on focus when all its light rays gathered by the camera hit the image plane at the same point

- In a pinhole device this happens to all scene points because of the very small size of the hole, so that the camera features an infinite Depth of Field (DOF)



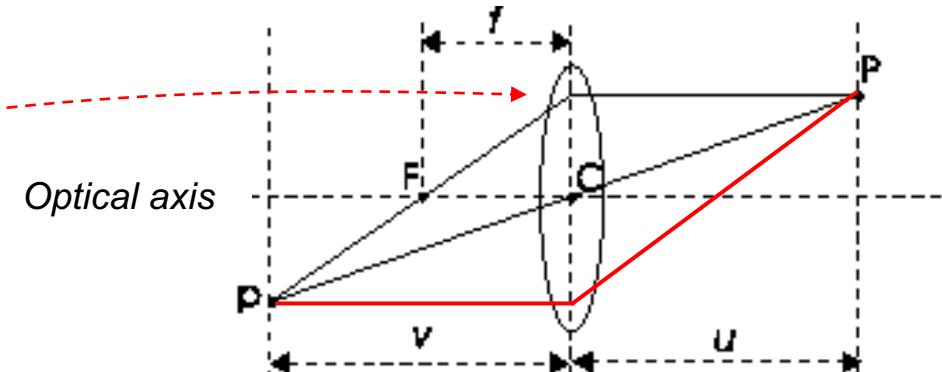
- The drawback is that such a small aperture allows gathering a very limited amount of light
- If a point is projected onto a circle instead of a point (bigger pinhole) the image is not sharp (not on focus) => you want to map all scene points to all image points
- Getting sufficiently bright images requires very long exposure times
  - If we cannot gather through aperture we have to gather (integrate) through time...
  - Long time => moving object? => only static scenes can be acquired by a pinhole device to avoid motion blur

# Lenses

- Use lenses to gather more light from a scene point and focus it on a single image point
  - This enables much smaller exposure times  
=> avoid motion blur in dynamic scenes
  - DOF is no longer infinite => only points across a limited range of distances can be simultaneously in focus in a given image
  - Cameras often feature complex optical systems, comprising multiple lenses
  - We will consider the approximate model known as **thin lens** equation:
- COLLECTS LIGHTS*
- LENS*
- CONCENTRATES LIGHTS*
- P*
- LENS REDUCE EXPOSURE TIMES AND DEPTH OF FOCUS*
- THIN LENSES MODEL*
- $$\frac{1}{u} + \frac{1}{v} = \frac{1}{f}$$

# Lenses

- To graphically determine the position of a focused image point we can leverage on two properties of thin lenses:
  - Rays parallel to the optical axis are deflected to pass through  $F$
  - Rays through  $C$  are undeflected
- If the image is on focus, the image formation process obeys to the perspective projection model:
  - the centre of the lens is the optical centre
  - the distance  $v$  acts as the effective focal length of the projection
- By fixing the position of the image plane...



$$\frac{1}{u} + \frac{1}{v} = \frac{1}{f}$$

$P$  : scene point

$p$  : corresponding focused image point

$u$  : distance from  $P$  to the lens

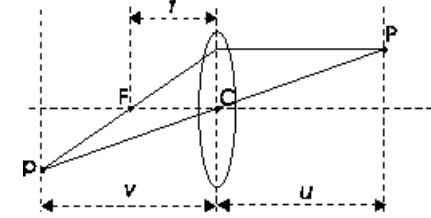
$v$  : distance from  $p$  to the lens

FixEΔ  $f$  : focal length (parameter of the lens)

$C$  : centre of the lens

$F$  : focal point (or focus) of the lens

# Lenses



- On one hand: choosing the distance of the image plane determines the distance at which scene points appear on focus in the image

$$\frac{1}{u} + \frac{1}{v} = \frac{1}{f} \rightarrow u = \frac{vf}{v-f}$$

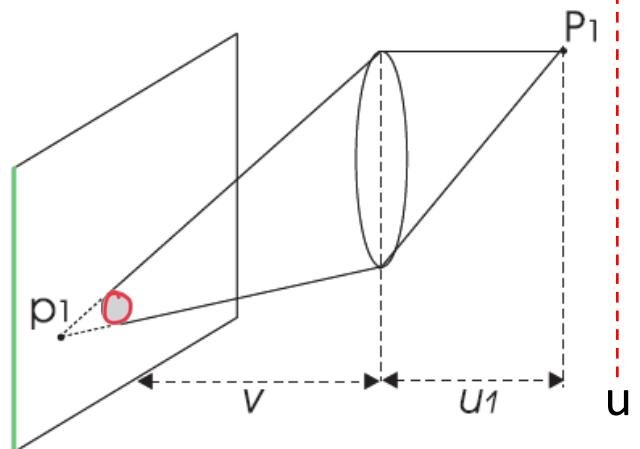
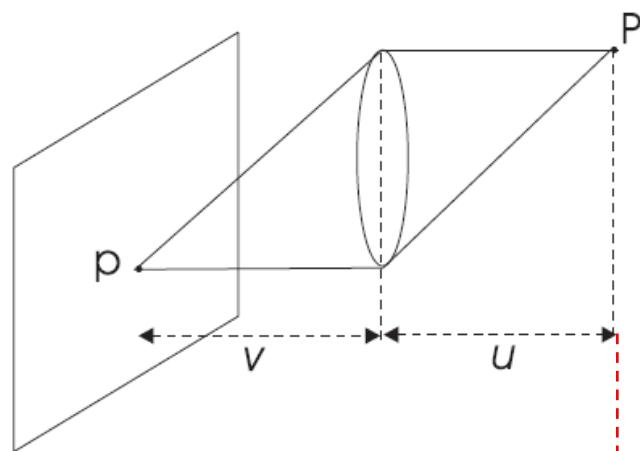
- On the other hand: to acquire scene points at a certain distance we must set the position of the image plane accordingly:

$$\frac{1}{u} + \frac{1}{v} = \frac{1}{f} \rightarrow v = \frac{uf}{u-f}$$

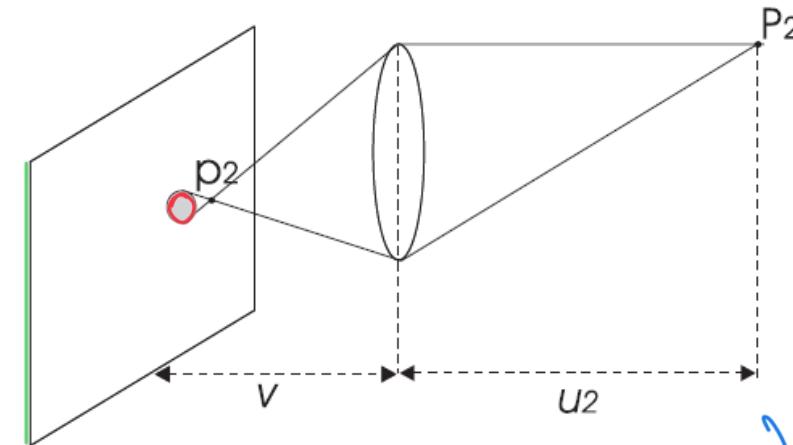
- Given the chosen position of the image plane, scene points both in front and behind the focusing plane will result out-of-focus, thereby appearing in the image as circles, known as **Circles of Confusion** or **Blur Circles**, rather than points
- The advantage of lenses is to have a small exposure time for capturing moving objects but we pay in terms of depth of field

# Lenses

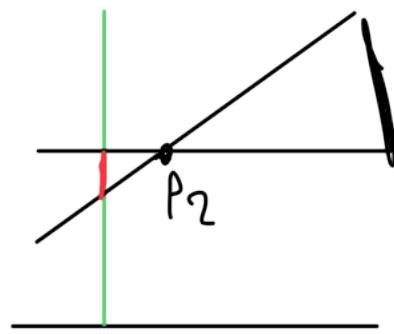
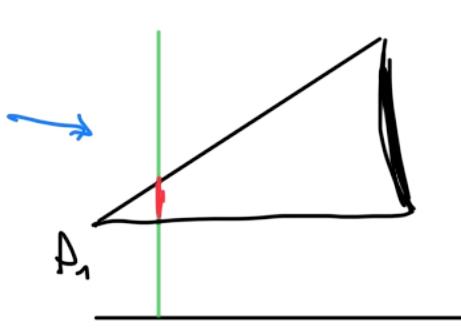
$P$  belongs to the focusing scene plane



$p_1$  is behind the image plane



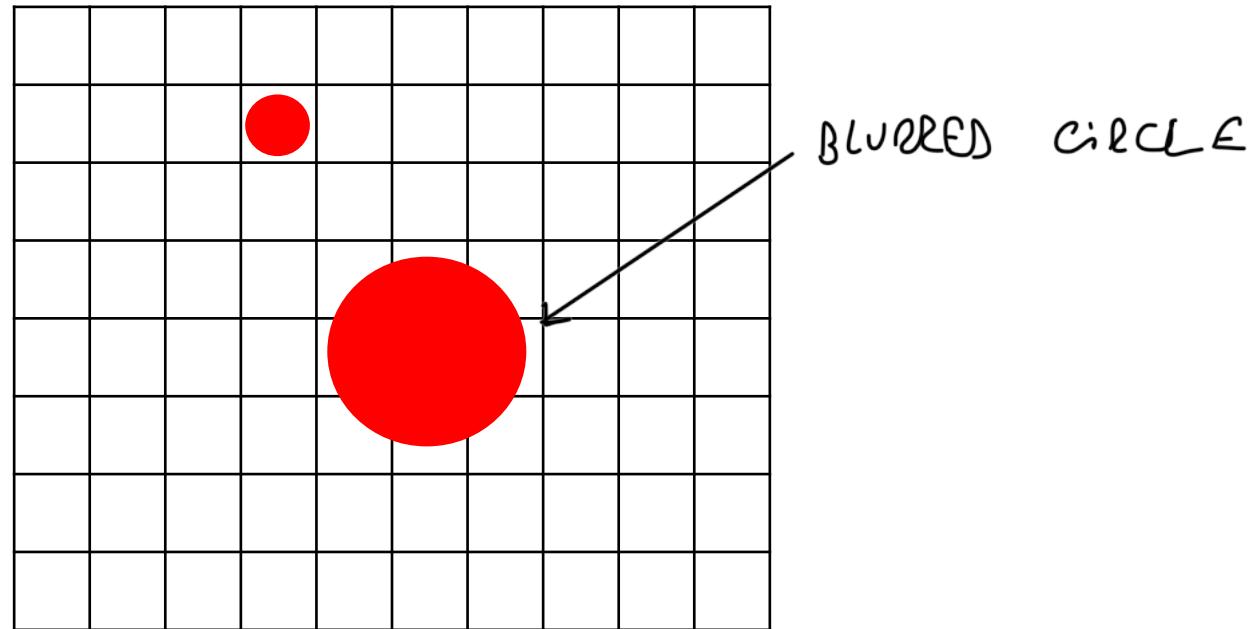
$p_2$  is in front of the image plane



# Diaphragm

---

- In theory, when imaging a scene through a thin lens, only the points at a certain distance can be on focus, all the others appear blurred into circles
  - However, as long as such circles are smaller than the size of the photosensing elements, the image will still look on-focus (i.e. the light is collected by a single pixel of the camera sensor)



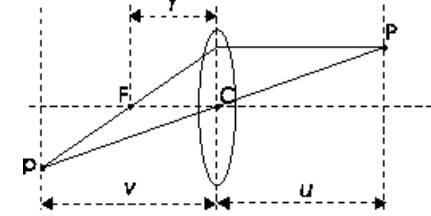
# Diaphragm

---

- In theory, when imaging a scene through a thin lens, only the points at a certain distance can be on focus, all the others appear blurred into circles
  - However, as long as such circles are smaller than the size of the photosensing elements, the image will still look on-focus (i.e. the light is collected by a single pixel of the camera CCD)
  - The range of distances across which the image appears on focus - due to blur circles being small enough - determines the DOF (Depth of Field) of the imaging apparatus
- Cameras often deploy an adjustable diaphragm (iris) to control the amount of light gathered through the *effective aperture* of the lens
  - Reduce aperture => less light => smaller blur circle
  - More aperture => more light => larger blur circle
- We close the diaphragm => increase depth of field => not enough light => increase exposure time => moving object => motion blur => still scenes

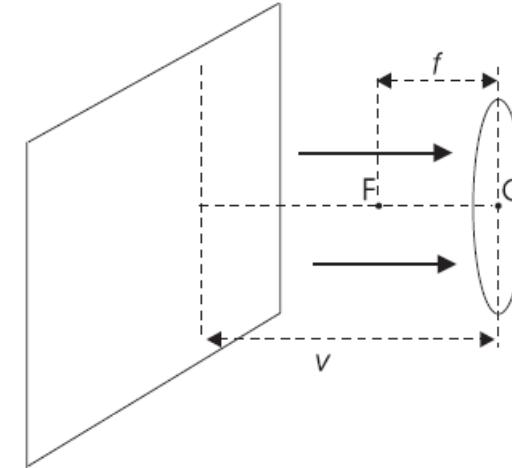
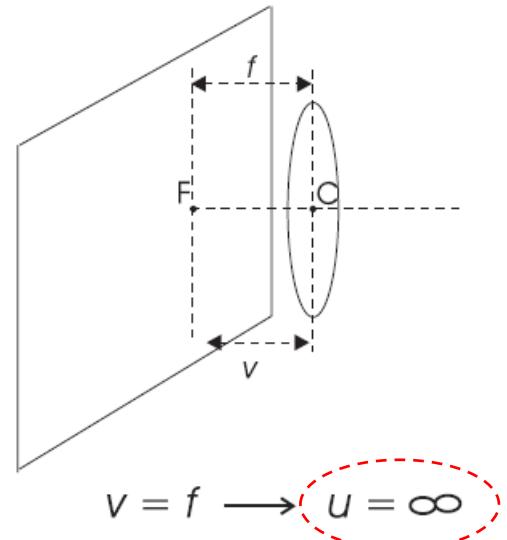


# Focusing Mechanism



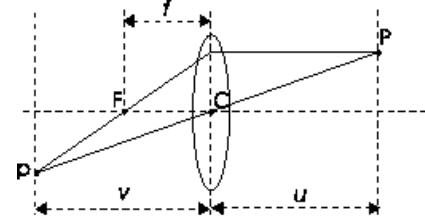
- To focus on objects at diverse distances:
  - mechanism that allows the lens (or lens subsystem) to translate along the optical axis with respect to the **fixed** position of the image plane

$$\cancel{\frac{1}{u} + \frac{1}{v} = \frac{1}{f}}$$



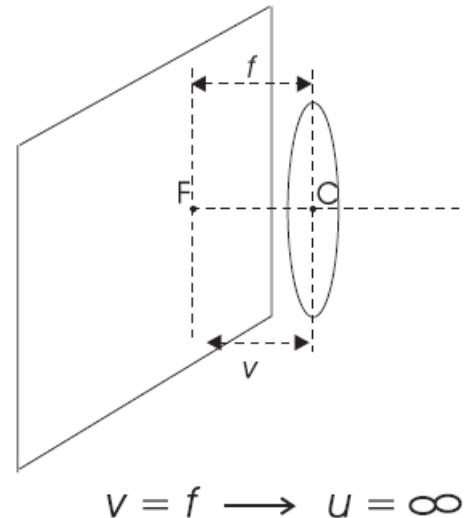
- At one end position ( $v=f$ ) the camera is focused at infinity

# Focusing Mechanism

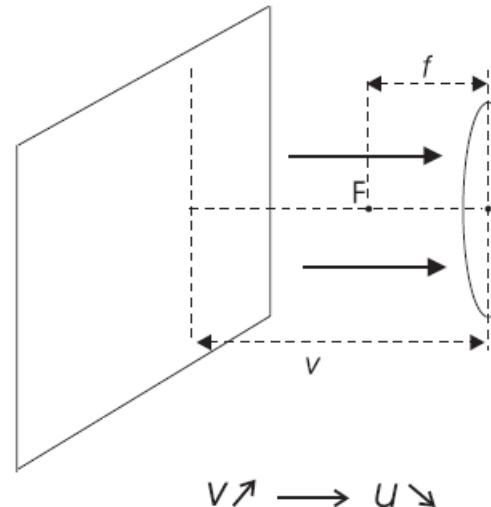


- To focus on objects at diverse distances:
  - mechanism that allows the lens (or lens subsystem) to translate along the optical axis with respect to the **fixed** position of the image plane

$$\frac{1}{u} + \frac{1}{v} = \frac{1}{f}$$



$$v = f \longrightarrow u = \infty$$



$$v > \longrightarrow u <$$

as long as we increase  
"v" we decrease "u"

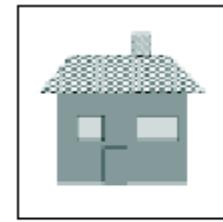
- At one end position ( $v=f$ ) the camera is focused at infinity
- The mechanism allows the lens to be translated farther away from the image plane up to a certain maximum value (the second end position), which determines the minimum focusing distance

# Image Digitization

- Generally speaking, the image plane of a camera consists of a planar sensor which converts the **irradiance** at any point into an electric quantity (e.g. a voltage)
  - Transduction from light to an electric quantity => "records" the light coming from a scene point

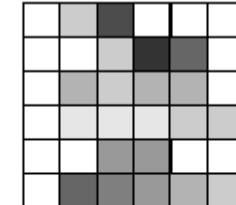
SENSORS SEE ONLY GRAY-SCALE

- How do we discretise it?



Continuous Image

BECAUSE OF FINITE  
NUMBER OF SENSORS  
Sampling



Sampled according to  
a two dimensional grid

Quantization

255	204	77	255	255	255
255	255	204	51	102	255
255	178	204	178	178	255
255	230	230	230	204	204
255	255	153	153	255	255
255	102	128	153	178	204

Sampled and  
Quantized Image

- Such a continuous "electric" image is sampled and quantized to end up with a digital image suitable to visualization and processing by a computer
  - The continuous voltages in the sampled image are gonna be quantized into a fixed number of levels

# Image Digitization

---

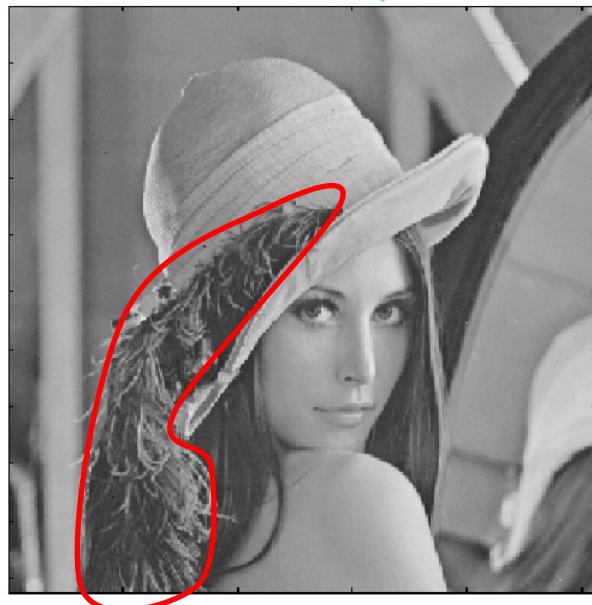
- **Sampling** – The planar continuous image is sampled along both the horizontal and vertical directions to pick up a 2D array (matrix) of  $N \times M$  samples known as pixels:

$$I(x, y) \implies \begin{bmatrix} I(0, 0) & I(0, 1) & \dots & I(0, M - 1) \\ \vdots & \vdots & & \vdots \\ I(N - 1, 0) & I(N - 1, 1) & \dots & I(N - 1, M - 1) \end{bmatrix}$$

- **Quantization** – The continuous range of values associated with pixels is quantized into  $l = 2^m$  discrete levels known as gray-levels
  - $m$  is the number of bits used to represent a pixel, with the memory occupancy (in bits) of a gray-scale image given by:  $B = N \times M \times m$ 
    - $m=8$  in gray-scale digital images, so that, e.g., a VGA format ( $480 \times 640$ ) image requires 300 Kbytes for storage while a 1mpx image requires 1 Mbytes.
    - colour digital images are instead typically represented within computers using 3 bytes per pixels (one byte for each of the RGB channels).
  - The more the better => more pixels + more bits per pixel => higher quality image

# Digitization vs. Image Quality

- The more bits we spend for its representation, the higher the quality of the digital image (we get a closer approximation to the ideal continuous image)
- This applies to both sampling as well as quantization parameters



512×512, l=256  
(original)



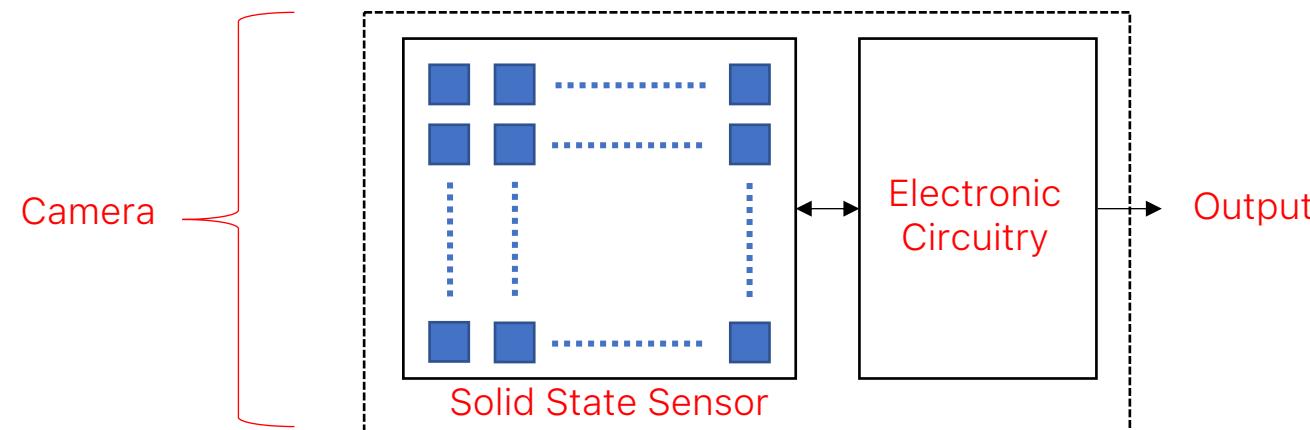
64×64, l=256  
(coarser sampling)



512×512, l=16  
(coarser quantization)

# Camera sensors

- The sensor is a 2D array of photodetectors (photogates or photodiodes)
  - During exposure time, each detector converts the incident light into a proportional electric charge (i.e. photons to electrons)
  - The companion circuitry reads-out the charge to generate the output signal, which can be either *digital* or analog
  - For digital: the camera includes also the necessary ADC circuitry

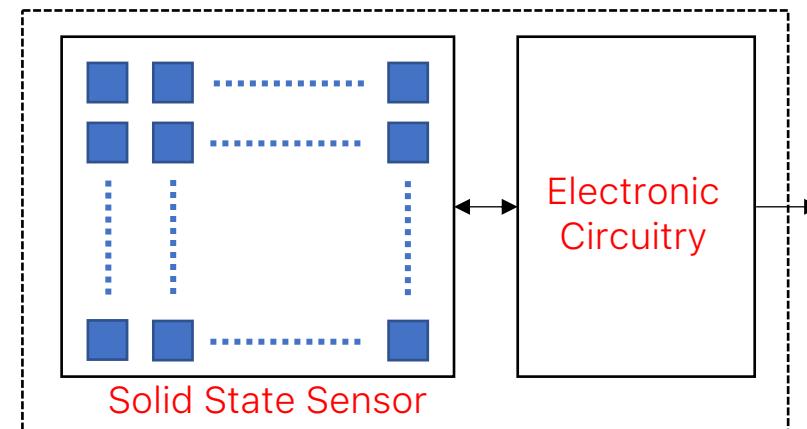


- Nowadays, there is never a continuous image in practice => the image is sensed directly as a sampled signal

# Camera sensors

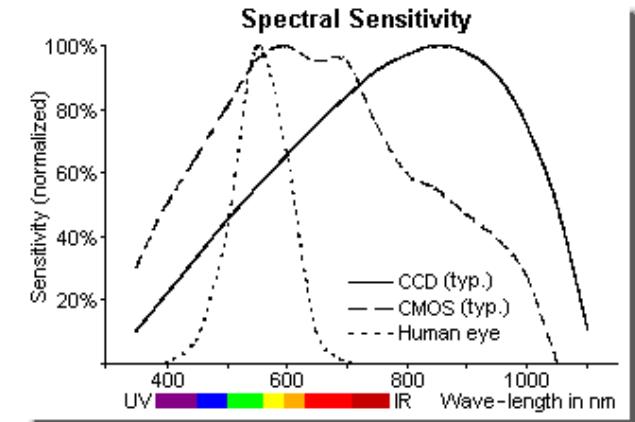
---

- Today, the two main sensor technologies are:
  - CCD (Charge Coupled Devices);
  - CMOS (Complementary Metal Oxide Semiconductor)
- Unlike CCD, CMOS technology allows the electronic circuitry to be integrated within the same chip as the sensor ("one chip camera")
  - This provides more compactness, less power consumption and often lower system cost.
- Unlike CCD, CMOS sensors allow an arbitrary window to be read-out without having to receive the full image
  - This can be useful to inspect or track at a higher speed a small Region Of Interest (ROI) within the image
- CCD technology typically provides higher Signal-to-Noise Ratio (SNR) , higher Dynamic Range (DR) and better uniformity



# Colour Sensors

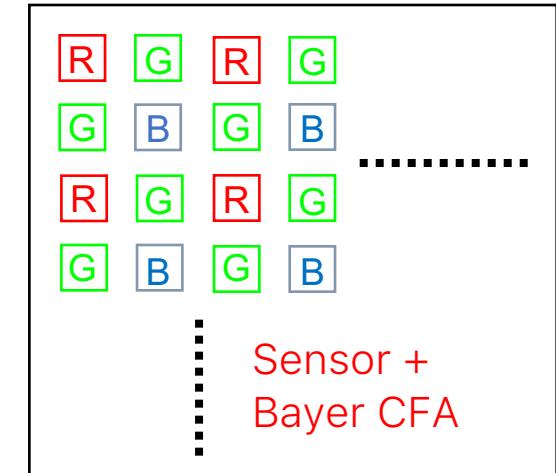
- CCD/CMOS sensors are sensitive to light ranging from near-ultraviolet (200 nm) through the visible spectrum (380-780 nm) up to the near infrared (1100 nm). The sensed intensity at a pixel results from integration over the range of wavelengths of the spectral distribution of the incoming light multiplied by the spectral response function of the sensor ==> **CCD/CMOS sensor cannot sense colour.**



- To create a colour sensor, an array of optical filters (Colour Filter Array) is placed in front of the photodetectors, so as to render each pixel sensitive to a specific range of wavelengths.

- In the most common, Bayer CFA, green filters are twice as much as red and blue ones to mimic the higher sensitivity of the human eye in the green range.
- To obtain an RGB triplet at each pixel, missing samples are interpolated from neighbouring pixels (demosaicking). However, the true resolution of the sensor is smaller due to the green channel being subsampled by a factor of 2, the blue and red ones by 4.

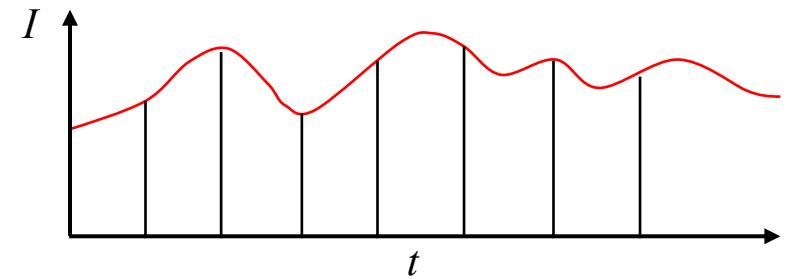
WE USE LESS PIXELS FOR RESOLUTION BECAUSE  
THEY BECOME SPECIFIC FOR THEIR ONLY  
WAVELENGTH AND LOSE THE OTHERS



- A "full resolution" – though more expensive - colour sensor can be achieved by deploying an optical prism to split the incoming light beam into 3 RGB beams sent to 3 distinct sensors equipped with corresponding filters.

# Signal to Noise Ratio

WE NEVER OBSERVE A CONSTANT SIGNAL



- *Signal-to-Noise Ratio (SNR)* – The intensity measured at a pixel under perfectly static conditions varies due to the presence of random noise (i.e. a pixel value is not deterministic but rather a random variable).
- The main noise sources are:
  - *Photon Shot Noise* – The time between photon arrivals at a pixel is governed by a Poisson statistic and thus the number of photons collected during exposure time is not constant.
  - *Electronic Circuitry Noise* – It is generated by the electronics which reads-out the charge and amplifies the resulting voltage signal.
  - *Quantization Noise* – related to the final ADC conversion (in digital cameras).  
→ LOST OF INFORMATION
  - *Dark Current Noise* – a random amount of charge due to thermal excitement is observed at each pixel even though the sensor is not exposed to light.
- The SNR can be thought of as quantifying the strength of the "true" signal with respect to the unwanted fluctuations induced by noise (i.e. the higher the better). It should be measured according to standard procedures and it is usually expressed either in *decibels* or *bits*:

$$\text{SNR}_{dB} = 20 \cdot \log_{10}(\text{SNR}); \quad \text{SNR}_{bit} = \log_2(\text{SNR})$$

# Dynamic Range

---

- Dynamic Range (DR) – If the sensed amount of light is too small, the “true” signal cannot be distinguished from noise
  - given  $E_{\min}$  - the minimum detectable irradiation (which depends on what?)
  - given  $E_{\max}$  - the saturation irradiation (i.e. the amount of light that would fill up the capacity of a photodetector).
- The DR of a sensor is defined as  $DR = \frac{E_{\max}}{E_{\min}}$ ; and, like the SNR, it is often specified in decibels or bits.
- As it is the case of the SNR, also for the DR the higher is the better. Indeed, the higher the DR the better is the ability of the sensor to simultaneously capture in one image both the dark and bright structures of the scene.
- An active research field in image processing deals with creating High Dynamic Range (HDR) images by combining together a sequence of images of the same subject taken under different exposure times (see e.g. <http://www.hdrsoft.com/index.html>).